

Part B

Part B.1 Formal Description

The models in part A suffer from one obvious problem: they fail at generalizing. One way to fix this is to decrease the noisiness of the data, and improve the quality of the local optima chosen during training.

To do this, we introduce Curriculum Learning (CL) into our models. We train our models first on easy data, and gradually increase in the difficulty. CL, as introduced by Wang et al. in A Survey on Curriculum Learning, is the combination of two algorithms: 1) a **difficulty measurer** that can sort the data from easiest to hardest, 2) a **training scheduler** that can choose when to increase the difficulty of the data the model is currently trained on.

The hope is that by choosing the right difficulty measurer and training scheduler, we can allow the training algorithms to capture the patterns in the data without being distracted by noise. The reason this works is that by denoising the data, we can find better local optima. The paper A Survey on Curriculum Learning uses figure 6 to explain this.

All of our added code can be found in the files: `difficulty_measurers.py`, `difficulty_measurers_dual_dual.py`, `model.py`, `main.py`.

Algorithm 1 Curriculum Learning Template

```
sortedData  $\leftarrow$  difficultyMeasurer(data)                                 $\triangleright$  Sorts the data from easiest to hardest
for N epochs do
    currentData  $\leftarrow$  TrainingScheduler(sortedData)                 $\triangleright$  Chooses the data to train on
    Model.train(currentData)
end for
```

Part B.1(a) Difficulty Measurers

The base difficulty measures are:

- **Number Of Question Entries:** We can make the assumption that the more questions a user has answered, the less noisy the data will be, and therefore the easier it is for the models to learn the pattern.
- **Autoencoder Bootstrapping:** The autoencoder we train in part a is capable of compressing and uncompressing the data. The key idea is that we can use this as a measure of noise. We simply sort the data in terms of how good the autoencoder is at reproducing it.
- **Question-Based Subject Correctness Entropy:** This difficulty measurer utilizes metadata "question_meta.csv". The intuition is that each question belongs to some subjects, and we rank subjects' difficulty (to our ML algorithm, not how conceptually hard a subject is!) by computing its entropy. We consider all questions that involve this subject, compute probability of correctness, then use that to compute entropy. After we derive subject entropy, each question will be ranked by the subject it belongs to that has the highest entropy. This is inspired by the "bucket effect".

In addition to the ones above, we also use a reverse ordering for purpose of testing, sorting the data from hardest to easiest, as well as a randomizer for control purposes.

Part B.1(b) Training Schedulers

The base training schedulers are:

Algorithm 2 Baby-steps Scheduler

```
sortedData  $\leftarrow$  difficultyMeasurer(data) ▷ Sorts the data from easiest to hardest
Divide sortedData into equal sized buckets
currentData  $\leftarrow \emptyset$ 
for each bucket in buckets do
    currentData.union(bucket)
    Model.train(currentData)
end for
```

- **Baby-Steps:** As described in Wang et al, the baby steps algorithm divides the sorted data into buckets, then after a fixed number of iterations, adds a new bucket to the current data.
- **Continuous Learning:** As described in Wang et al, instead of dividing the data into chunks, we can instead, on the i th epoch, train the model on the easiest $100\lambda(i)\%$ of the data. λ is a function of the number of epochs we have trained for, and is commonly known as the pacing function.

In our experimentation, we use the following two pacing functions. These are inspired by and modified from Wang et al. A linear pacing function:

$$\lambda_{\text{Linear}}(t) = \min(1, \lambda_0 + st) \quad (1)$$

where λ_0 is the initial proportion of the data, and s is the amount we increase every iteration.

A root pacing function:

$$\lambda_{\text{Geom}}(t) = \min(1, \sqrt{2st + \lambda_0}) \quad (2)$$

Where s is the rate of growth. This solves the differential equation $\frac{d\lambda}{dt} = \frac{s}{\lambda}$, so that the amount of new data we give is inversely proportional to the amount of data as the data gets harder.

- **Validation Learning:** We can use a pacing function that depends not on the epoch, but on the current accuracy on the validation set. Indeed,

$$\lambda_{VL}(a) = \min(1, a/g) \quad (3)$$

where a is the current validation accuracy, and g is a constant we called the goal accuracy.

Part B.2 Figure or Diagram

See figure 7 for a succinct description of the project.

Part B.3 Comparison and Description

Changing the order of the data clearly makes a difference, but it also clearly isn't a big difference.

In comparison with our base line models, we see that the test accuracy of IRT is actually improved by 0.2% to 0.7098504092576913 using the entropy difficulty measurer! However, the autoencoder seems to suffer from being fed the data in segments, decreasing to 0.6886. Regardless of the order of the data, we cannot get 0.6910 test accuracy in part a.

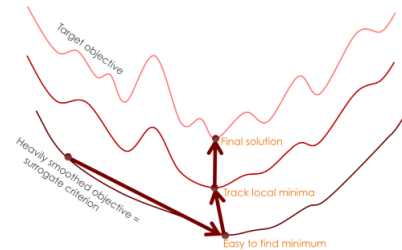


Figure 6: How CL improves the quality of local minima (Taken from Wang et al.)

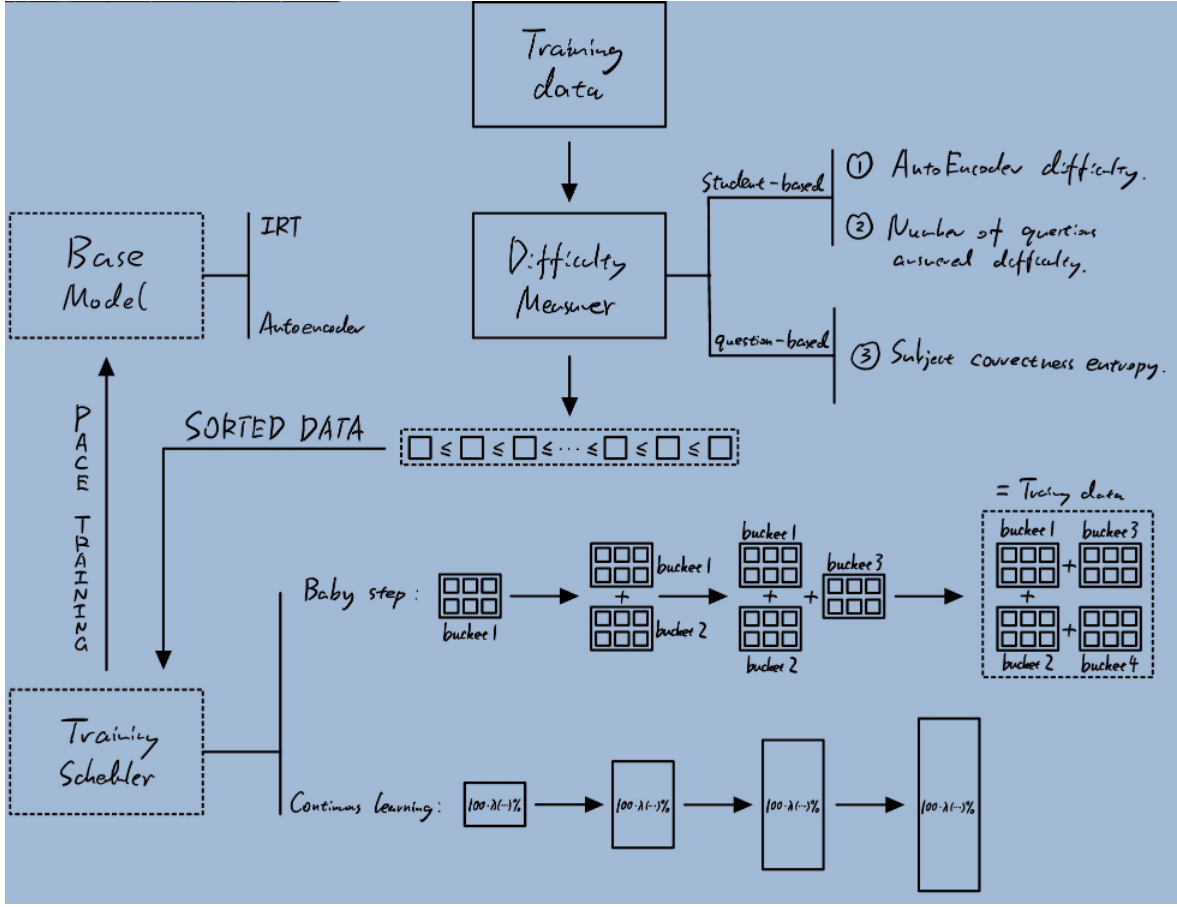


Figure 7: An overview of CL as we implemented it. Each box is a component, and the arrows between them represent the flow of data.

Training Scheduler	Test accuracy with autoencoder	Test accuracy with number of entries
Continuous Validation	0.6836014676827548	0.6796500141123342
Linear Validation	0.6824724809483489	0.6830369743155518
Baby Steps	0.6886819079875811	0.6802145074795372

Table 1: Autoencoder Results

Our first hypothesis is that using CL will guide the optimization algorithm to better local optima. To test this with a control, we ran experiments where we use CL with a difficulty measurer, without a difficulty measurer (i.e. randomly sorting the data), and with a reverse difficulty measurer.

For an autoencoder with 15 buckets and 40 iterations per bucket with baby steps: using the autoencoder difficulty measure, we get a test accuracy: 0.6802145074795372. Sorting from hardest to easiest results in a test accuracy of 0.6700536268698842. Randomly sorting the data gives us a test accuracy of 0.6740050804403048. See figures 8 to 10 for the results. CL is better than random, especially early on, but the difference is small at the end. The only difference is the order, suggesting that with a good difficulty measure, we indeed find better local optima.

The next hypothesis is that using an autoencoder to sort the data is a good idea, that this sort of autoencoder bootstrapping is useful. Indeed, comparing the test accuracy, we do indeed see this behaviour, with the auto encoder difficulty measurer performing consistently better than the number of entries difficulty measure in the test accuracy in table 1, as well as doing better than random.

We also see this in the qualitative behaviour in figures 8 to 10. Interestingly, we see that training the

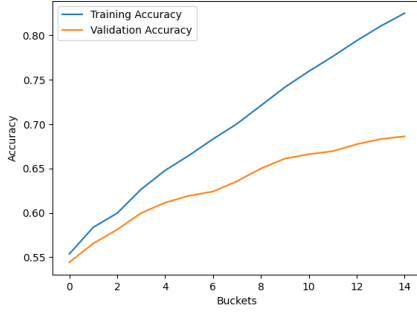


Figure 8: Autoencoder with Autoencoder bootstrapping

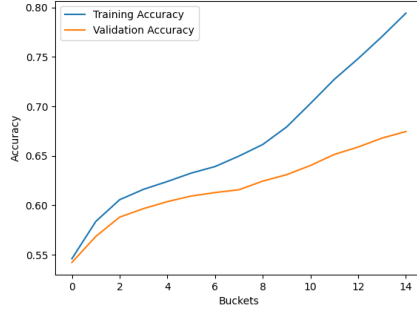


Figure 9: Autoencoder with random sorting

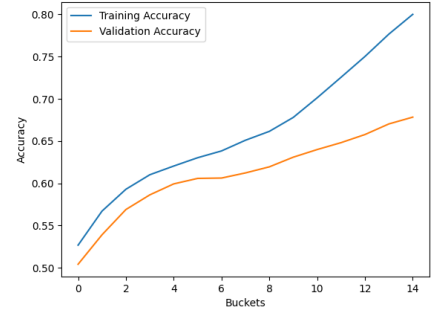


Figure 10: Autoencoder with Autoencoder bootstrapping backwards

autoencoder with randomly sorted data, or with data sorted from worst performance on the autoencoder to the best performance has a plateau in the middle. But with autoencoder bootstrapping there seems to be continuous, faster improvement. Sorting in an unhelpful way withholds useful information from the model, providing evidence that the difficulty measurer does indeed help our model learn.

We have also investigated the effect of CL on IRT model. We see that autoencoder difficulty is the worst here, and only the entropy difficulty measure is consistently better than random. The changes are very small, IRT seems to be more or less robust in terms of the ordering of data.

Part B.4 Limitations

1. **Useless on convex models:** CL is only meaningful when working with models that are non-convex. For example, there is no point of this scheme linear regression as it wastes time and computing resources while linear regression alone could just find the global minimum.
2. **Hard to choose scheduler and measurer:** The difficulty measure is more heuristic than mathematical. During our testing, we are often surprised by the test accuracy when trying out different options of implementing CL algorithm. The difficulty measures that made sense to us often yield terrible performance. As explained in the paper, this is due to the decision boundary of human are different from that of computer, so what is intuitive to us isn't as helpful to our ML algorithm. This means that we have to try out many difficulty measures then pick one that works the best, even sometimes against our intuition.

One open problem is to figure out which measurers are best for each autoencoder and IRT.

3. **Data Sparsity and Faulty Generalization:** CL could be misleading to the base model, especially with sparse data, because it withholds information. For instance, in baby steps we only train on 113 data points in the first bucket. There may be questions that are not known to be answered by all students in that bucket, or are answered incorrectly by all students in that bucket (i.e. “zero column” in the matrix). This may of mislead the model into thinking this question is impossible, which may explain why the autoencoder does worse with segmented data.
4. **Time:** It takes so much longer to train with CL. There is much more work to do, including having to filter through the training data at each epoch, and restart training from the beginning. It is much harder to experiment this way, especially with the autoencoder.

Possible extensions include applying curriculum learning to an ensemble, building autoencoder bootstrappers that by encode questions to order them by noisiness, trying a larger variety of pacing functions, and gathering more data to deal with data sparsity.